# DIAG: Data Intensive Academic Grid

## A computational platform for bioinformatics analyses and training

Owen White

Anup Mahurkar

- Who we are

- Motivation for MRI-R$^2$

- System Description

- Highlight of users and applications

- Timeline

- Challenges

# *Who we are*

- Part of the University of Maryland School of Medicine in Baltimore
- Quasi-independent institute founded about 2 ½ years ago
- IGS has ~20 faculty members
- Total size of ~100
- Small to medium size sequencing center
- Mix of sequencing platforms including Sanger, Illumina, and Roche 454
- Areas of focus: microbial genomics, human genomics, and metagenomics
- http://www.igs.umaryland.edu

**INSTITUTE FOR GENOME SCIENCES**

- Sequencing technology improvements
  - 1st Generation
    - Sanger-based capillary sequencing with throughput of ~500 Kb/run
  - Lots of players on the horizon
    - 454, Illumina, Solid, Helicos, Visigen, PacBio, Complete Genomics, Oxford Molecular, NABsys, IBM, Life Technologies
  - 2nd Generation
    - 454 pyrosequencing – 160-320 Mb/run (8 hour run)
    - Illumina HiSeq2000 200 Gb/run (8 day run) or 25 Gb/day
    - ABI SOLiD 200-300 Gb/run
  - 3rd Generation
    - Pacific Biosciences,
      - Human genome < 15 mins for $100 by 2014
    - Complete Genomics
      - 50 human genomes completed
      - 500 on order to be completed in 2010
      - 18 genomes per 11 day run will rise to 120 genomes per run
      - Roadmap calls for 1 million genomes/year capacity in next few years

**INSTITUTE FOR GENOME SCIENCES**

- Democratization of sequencing with 100s if not 100s of these sequencers being sold

- Burgeoning data sets
  - 1000 human genomes data
  - 1000s of cancer genomes
  - 2000 bacterial genomes
  - 10s of plant genomes
  - Metagenomes for various environments including soil, ocean, air, human body, etc

**INSTITUTE FOR GENOME SCIENCES**

# *Motivation*

- Newer applications/projects
  - 1000 genomes project
  - Cancer Genomics
  - Transcriptomics
  - Epigenetics
  - Population genomics
  - Metegenomics
    - 16S rRNA based community classification
    - Whole genome sequencing of metagenomes
    - Metatranscriptomics
- Computationally intensive
  - A *de novo* assembly will take order of a day or two on a 128 node cluster using tools like *ABySS*
  - Metagenomic annotation using Blast, and HMM search will take over 9000 hours on a single core for data generated by a 454 in 8 hours
- Lack of proximity between reference data sets and computational resources

**INSTITUTE FOR GENOME SCIENCES**

# *DIAG Highlights*

- Over 20 users from 13 US and international institutions
- Diverse applications
  - Microbial genomics
    - Annotation
    - Comparative Genomics
  - Plant Genomics
    - Assembly
    - Annotation
    - Transcriptomics
  - Metagenomics
    - Marine
    - Human
    - Plant
    - Environment
  - Proteomics
  - Livestock Research

**INSTITUTE FOR GENOME SCIENCES**

- Sources include: GenBank, RefSeq, Uniprot (including Swiss-Prot), UniProtKB, PDB, PIR, Ensembl, EMBL, CAMERA, MG-RAST, RAST sub-systems and Greengenes.
- Non-redundant protein and nucleotide data sets
  - Custom data sets for bacteria, viruses, eukaryotes, mammals generated by PANDA
- Data from 1000 genomes project
- Metagenomics data sets
  - CAMERA
  - Virome
  - Human Microbiome Project
- Transcriptomes for various model organisms including human, mouse, drosophila, Arabidopsis, etc.

## *System Highlights*

- 100-125 high-throughput compute nodes
  - 400 GB local storage per node
  - 48 GB RAM per node
  - Intel/AMD processors
- 5 high-performance compute nodes
  - 1 TB local storage per node
  - 12-16 cores per node
  - ~12-18 GB RAM per core
  - InfiniBand QDR interconnect
- 400-600 Tera Bytes (TB) storage
  - High-performance, grid-attached parallel file system (GPFS, Lustre, Isilon, Panasas)
  - Possibly hierarchical storage
  - Archival upon request

**INSTITUTE FOR GENOME SCIENCES**

# System Highlights

# TIGR Workflow

- Open source workflow management software written in Java at
- Written at The Institute for Genomic Research/J Craig Venter Institute
- Currently maintained by our group at IGS
- http://tigr-workflow.sourceforge.net/
- In use for over 8 years
- A GUI for authoring pipelines
- Multi-threaded workflow execution engine
- A GUI for monitoring pipeline execution
- Job execution on a single host or on a grid
- Built in support for Condor and Sun Grid Engine
- Working on support for PBS and a generic grid through DRMAA

**INSTITUTE FOR GENOME SCIENCES**

- Open source web-based pipeline creation and monitoring platform
- http://sourceforge.net/projects/ergatis/
- Uses TIGR Workflow as job execution engine
- Generates workflow XML directly
- Has over 50 commonly used bioinformatics applications used for genome assembly, sequence searches
- In use for over 8 years at a number of institutions
- Used in a number of major research projects including the analysis of Annotation of Plasmodium Genomes, Global Ocean Survey Metagenomics

- A genomics tool for automated and portable sequence analysis using Virtual Machines and Cloud computing
- PI: Florian Fricke
- Tools for microbial assembly, annotation, comparative analysis
- Tools for microbial and viral metagenomic community profiling
- Tools for whole metagenome analysis and annotation
- http://clovr.igs.umaryland.edu

# Selects Users and Applications

| PI | Institution | Application |
|---|---|---|
| Fricke | IGS | VM based pipelines |
| Giglio | IGS | Microbial Annotation |
| Silva | IGS | Evolutionary Analysis |
| Meyer | Univ. Chicago/ANL | Metagenomics |
| Andersen | LBNL | Metagenomics |
| Wommack | Univ. Delaware | Viral Metagenomics |
| Buell | Michigan State | Plant Genomics |
| Tettelin | IGS | Comparative Pangenome Analysis |

**INSTITUTE FOR GENOME SCIENCES**

# *Select Contributors of Tools/Pipelines/Data*

| PI | Institution | Contribution |
| --- | --- | --- |
| Fricke | IGS | Analysis VMs (CloVR) |
| Wommack | Univ. Delaware | Virome |
| Pearson | Univ. Virginia | FASTA |
| Tettelin | IGS | Sybil |
| Giglio | IGS | Microbial Annotation Pipeline |
| Field | | Bio-Linux |
| Ellisman | Univ. California | CAMERA data sets, API |
| Meyer | Univ. Chicago/ANL | MG-RAST |
| Andersen | LBNL | Metagenomics classification tools |

**INSTITUTE FOR GENOME SCIENCES**

- Prototype system in Spring 2010
  - ~100 cores, 10TB storage
  - Direct access
  - Ergatis
  - Nimbus
  - OSG compute element
- Hardware Acquisition Fall 2010
- Deployment and testing Winter 2010
- Fully operational Early 2011

INSTITUTE FOR **GENOME SCIENCES**

- Deploying the entire cluster as a Cloud computing platform
- Authentication and authorization for multi-institution user base
- Dynamic scheduling of various access methods
- Identifying optimal data sets
- Bandwidth from institutions to upload data
- Space management

**INSTITUTE FOR GENOME SCIENCES**

- All the users and contributors
- PI and Co-PIs
  - Owen White (PI)
  - Eric Wommack
  - Sam Angiuoli
  - Jacques Ravel
- Institute for Genome Sciences
  - Engineering Team
    - Victor Felix
    - Joshua Orvis
    - Sam Angiuoli
  - IT Team
    - Dave Kemeza
    - Brian Cotton
    - Daniel Strassler
- University of Maryland School of Medicine
    - James McNamee
    - Scott Hunsinger

- University of Maryland Institute for Advanced Computer Studies
  - Fritz McCall
- RENCI
  - John McGee
- University of Chicago/Argonne National Lab
  - Kate Keahey

# Funding provided by: